# dsMTL - a computational framework for privacy-preserving, distributed multi-task machine learning

Han Cao[1#], Youcheng Zhang[2#], [*alphabetical order starts*] Jan Baumbach[3,4], Paul R Burton[5], Dominic Dwyer[6], Nikolaos Koutsouleris[6], Julian Matschinske[3], Yannick Marcon[7], Sivanesan Rajan[1], Thilo Rieg[1], Patricia Ryser-Welch[5], Julian Späth[3], [*alphabetical order ends*], The COMMITMENT consortium, Carl Herrmann[2*], Emanuel Schwarz[1*]

[1] Department of Psychiatry and Psychotherapy, Central Institute of Mental Health, Medical Faculty Mannheim, Heidelberg University, Mannheim, Germany.

[2] Health Data Science Unit, Medical Faculty Heidelberg & BioQuant, Heidelberg, 69120, Germany.

[3] Chair of Computational Systems Biology, University of Hamburg, Hamburg, Germany

[4] Computational Biomedicine Lab, Department of Mathematics and Computer Science, University of Southern Denmark, Odense, Denmark

[5] Population Health Sciences Institute, Newcastle University, Newcastle upon Tyne, United Kingdom

[6] Department of Psychiatry and Psychotherapy, Section for Neurodiagnostic Applications, Ludwig-Maximilian University, Munich 80638, Germany.

[7] Epigeny, St Ouen, France

*To whom correspondence should be addressed:    emanuel.schwarz@zi-mannheim.de

carl.herrmann@bioquant.uni-heidelberg.de

[#] These authors contributed equally to the study

**Consortium authors**

**The COMMITMENT consortium:**

Emanuel Schwarz[1*], [*alphabetical order starts*] Dag Alnæs[2], Ole A. Andreassen[2], Han Cao[1], Junfang Chen[1], Franziska Degenhardt[3,4], Daria Doncevic[5], Dominic Dwyer[6], Roland Eils[5,7], Jeanette Erdmann[8], Carl Herrmann[5], Martin Hofmann-Apitius[9], Nikolaos Koutsouleris[6,10], Alpha T. Kodamullil[9], Adyasha Khuntia[6], Sören Mucha[8], Markus M. Nöthen[3,11], Riya Paul[6], Mads L. Pedersen[12], , Heribert Schunkert[13], Heike Tost[1], Lars T. Westlye[2,12], Youcheng Zhang[5], [*alphabetical order ends*], Andreas Meyer-Lindenberg[1*]

[1] Department of Psychiatry and Psychotherapy, Central Institute of Mental Health, Medical Faculty Mannheim, Heidelberg University, Mannheim, Germany.

[2] Norwegian Centre for Mental Disorders Research (NORMENT), Division of Mental Health and Addiction, Oslo University Hospital and Institute of Clinical Medicine, University of Oslo, Oslo, Norway.

[3] Institute of Human Genetics, University of Bonn, School of Medicine & University Hospital Bonn, Bonn, Germany

[4] Department of Child and Adolescent Psychiatry, Psychosomatics and Psychotherapy, University Hospital Essen, University of Duisburg-Essen, Duisburg, Germany

[5] Health Data Science Unit, Medical Faculty Heidelberg and BioQuant, Heidelberg, 69120, Germany.

[6] Department of Psychiatry and Psychotherapy, Section for Neurodiagnostic Applications, Ludwig-Maximilian University, Munich 80638, Germany.

[7] Center for Digital Health, Berlin Institute of Health and Charité, Berlin, 10117, Germany.

[8] Institute for Cardiogenetics, University of Lübeck, DZHK (German Research Centre for Cardiovascular Research), partner site Hamburg/Lübeck/Kiel, and University Heart Center Lübeck, Lübeck, Germany.

[9] Fraunhofer Institute for Algorithms and Scientific Computing (SCAI), Sankt Augustin, 53754, Germany.

[10] Max-Planck Institute of Psychiatry, Munich, Germany

[11] Department of Genomics, Life & Brain Center, University of Bonn, Bonn, Germany.

[12] Department of Psychology, University of Oslo, Oslo, Norway

[13] Department of Cardiology, Deutsches Herzzentrum München, Technische Universität München, Munich Heart Alliance (DZHK), Germany.

57    **Abstract**

58    Multitask learning allows the simultaneous learning of multiple 'communicating' algorithms. It is

59    increasingly adopted for biomedical applications, such as the modeling of disease progression. As data

60    protection regulations limit data sharing for such analyses, an implementation of multitask learning on

61    geographically distributed data sources would be highly desirable. Here, we describe the development

62    of dsMTL, a computational framework for privacy-preserving, distributed multi-task machine learning

63    that includes three supervised and one unsupervised algorithms. dsMTL is implemented as a library

64    for the R programming language and builds on the DataSHIELD platform that supports the federated

65    analysis of sensitive individual-level data. We provide a comparative evaluation of dsMTL for the

66    identification of biological signatures in distributed datasets using two case studies, and evaluate the

67    computational performance of the supervised and unsupervised algorithms. dsMTL provides an easy-

68    to-use framework for privacy-preserving, federated analysis of geographically distributed datasets,

69    and has several application areas, including comorbidity modeling and translational research focused

70    on the simultaneous prediction of different outcomes across datasets. dsMTL is available at

71    https://github.com/transbioZI/dsMTLBase                (server-side            package)            and

72    https://github.com/transbioZI/dsMTLClient (client-side package).

73

74

75

76

77

78

79

80   **Introduction**

81   The biology of many human illnesses is encoded in a vast number of genetic, epigenetic, molecular,

82   and cellular parameters. The ability of Machine Learning (ML) to jointly analyze such parameters and

83   derive algorithms with potential clinical utility has fueled a massive interest in biomedical ML

84   applications. One of the fundamental requirements for such ML algorithms to perform well is the

85   availability of data at a large scale, a challenge of steadily declining importance due to the ever-

86   increasing availability of biological data[1-3]. As data can often not be freely exchanged across institutions

87   due tothe need for protection of the individual privacy, the utility of 'bringing the algorithm to the data'

88   is becoming apparent. Technological solutions for this task have thus risen in popularity and exist in

89   various forms. One of the most straightforward approaches is the so-called federated ML, where

90   algorithms are simultaneously learned at different institutions and optimized through a privacy-

91   preserving exchange of parameters. Other approaches for this task include the training of ML

92   algorithms on temporarily combined data stored in working memory[4] or the more recently introduced

93   'swarm-learning' approach[5]. One commonality of most ML algorithms, federated or not, is the

94   assumption that all investigated observations (e.g. illness-affected individuals) represent the same

95   underlying population. However, in biomedicine, this is rarely the case, as biological and technological

96   factors frequently induce cohort-specific effects that limit the ability to identify reproducible biological

97   findings. Multitask Learning (MTL) can address this issue through the simultaneous learning of

98   outcome (e.g. diagnosis) associated patterns across datasets with dataset-specific, as well as shared,

99   effects. Multi-task learning has numerous exciting application areas, such as comorbidity modeling,

100  and has already been applied successfully for e.g. disease progression analysis[6].

101  Here, we describe the development of dsMTL ('Federated Multi-Task Learning for DataSHIELD'), a

102  package of the statistical software R, for **Fe**derated **M**ulti-**T**ask **L**earning (FeMTL) analysis (**Figure 1**) .

103  dsMTL was developed for DataSHIELD[7], a platform supporting the federated analysis of sensitive

104  individual-level data that remains stored behind the data owner's firewall throughout analysis[8]. dsMTL

105  includes three supervised and one unsupervised federated multi-task learning algorithms that extend

106     algorithms previously developed for non-federated analysis (for R implementations, see [9,10]).

107     Specifically, the **dsMTL_L21** approach allows for cross-task regularization, building on the popular

108     LASSO method, in order to identify outcome-associated signatures with a reduced number of features

109     shared across tasks. The non-federated version of this approach has previously been applied to

110     simultaneously predict multiple oncological outcomes using gene expression data[11]. The **dsMTL_trace**

111     approach constrains the coefficient vectors in a low-dimensional space during the training procedure

112     to penalize the complexity of task relationships, resulting in an improved generalizability of the models.

113     In a non-federated implementation, this method has previously been used to predict the response to

114     different drugs, and the identified models showed a high degree of interpretability in the context of

115     the represented drug mechanism[12]. **dsMTL_net** incorporates the task relationships that can be

116     described as a graph, in order to improve biological interpretability. In a non-federated version, this

117     technique has previously been used for the integrative analysis of heterogeneous cohorts[13] and for the

118     prediction of disease progression[14]. The **dsMTL_iNMF** approach is an unsupervised, integrative non-

119     negative matrix factorization method that aims at factorizing the cohorts' data matrices into shared

120     and dataset-specific components. Such modeling has been applied to explore dependencies in multi-

121     omics data for biomarker identification[10,15]. In addition to the FeMTL methods, we also implemented

122     a federated version of conventional Lasso (dsLasso) [16] in dsMTL package due to its wide usage in

123     biomedicine and as a benchmark for testing the performance of the federated MTL algorithms.

124     To explore the utility of the dsMTL algorithms, we used a network comprising three servers. These

125     servers hosted simulated data with variable degrees of cross-dataset heterogeneity, in order to test

126     the ability of the MTL algorithms to suitably characterize shared and specific biological signatures. In

127     addition, we analyzed actual RNA sequencing and microarray data across the three-server network, to

128     show that the accurate analysis can be performed in acceptable runtime using dsMTL in real network

129     latency.

130

**Results**

Here we show the results for two case studies. The first case study aims at demonstrating the utility of the supervised **dsMTL_L21** algorithm to identify 'heterogeneous' target signatures across the data network. With 'heterogeneous' we describe signatures that involve the same features (e.g. genes) but with potentially differing signs (indicating differential directions of influences) across datasets. In contrast, 'homogeneous' signatures relate to the same features and signs across datasets. The second case study focuses on the unsupervised **dsMTL_iNMF** method and explores the utility of the federated implementation, compared to the aggregation of local NMF models, to disentangle shared and cohort-specific components across datasets. For all case studies, we evaluated the signature identification accuracy as the major metric. For predictions of clinical outcomes, the prediction accuracy was also demonstrated.

*Case study 1 – distributed MTL for identification of heterogeneous target signatures*

With the aim to identify 'heterogeneous' signatures, we compared the performance of dsMTL_L21, dsLasso and the bagging of glmnet models. As part of this, we explored the sensitivity of these methods to different sample sizes ($n$) relative to the gene number ($p$). **Figure 2** shows the resulting prediction performance and gene selection accuracy, each averaged over 100 repetitions. dsLasso showed the worst prediction performance in this heterogeneous setting, and dsMTL_L21 slightly outperformed the aggregation of local models (glmnet). Similarly, the gene selection accuracy of dsLasso was inferior to that of dsMTL_L21 and glmnet-bagging, which showed similar performance when the sample size is sufficiently large, e.g. the number of subjects approximately equal to the number of genes (n/p ~1). However, with a decreasing n/p ratio, dsMTL_L21 showed an increasing superiority over the other methods, especially for n/p=0.15, where the gene selection accuracy of dsMTL_L21 was over 2.8 times higher than that of the bagging technique.

156    *Case study 2 – distributed iNMF for disentangling shared and cohort-specific signatures*

157    **Figure 3** shows the performance of distributed and aggregated local NMF methods for disentangling

158    shared and cohort-specific signatures from multi-cohort data, given different 'severities' of the

159    signature heterogeneity. For both types of signatures, dsMTL_iNMF outperformed the ensemble of

160    local NMF models for any heterogeneity severity setting. Notably, even with increasing heterogeneity,

161    the accuracy of dsMTL_iNMF to capture shared genes remained stable at approximately 100%,

162    illustrating the robustness of dsMTL_iNMF against the heterogeneity's severity shown in **Figure 3c**. In

163    contrast, for the ensemble of local NMF, the gene selection accuracy of the shared signature

164    continuously decreased to approximately 50% (20% of outcome-associated genes were shared among

165    cohorts), while the gene selection accuracy of cohort-specific signatures continuously increased to 75%

166    (20% of outcome-associated genes were shared among cohorts ) as shown in **Figures 3a** and **3b**.

167

168    *Efficiencyof supervised dsMTL*

169    We aimed at determining the efficiency of supervised dsMTL using the real molecular data and the

170    actual latency of a distributed network. Using a three-server scenario (see **Table 2 Supplementary**

171    **Results**; two servers at the Central Institute of Mental Health, Mannheim; one server at BioQuant,

172    Heidelberg University) we analyzed four case-control gene expression datasets of patients with

173    schizophrenia and controls (median n=80; 8013 genes). **Supplementary Table 3** shows the comparison

174    between dsLasso and mean-regularized dsMTL_net, which were trained (cross-validation + training)

175    and tested in approximately 8min and 10min, respectively, with the time-difference being due to the

176    increased network access of dsMTL. The prediction accuracy of dsMTL was slightly higher than that of

177    dsLasso, consistent with our previous study[13]. Regarding model interpretability, dsLasso captured a

178    signature comprising 38 genes but could not distinguish shared and cohort-specific effects. Mean

179    regularized dsMTL identified a signature with 10 genes shared among all cohorts, with 163 genes

180    shared by two cohorts, as well as three cohort-specific signatures comprising 1532 genes.

181

182    *Efficiency of unsupervised dsMTL*

183    The cohorts and server information is shown in **Supplementary Table 4**. It took 34.9 minutes (1,003

184    times network accesses) to train a dsMTL_iNMF model with 5 random initializations (~7 min for each

185    initialization). The factorization rank k=4 was selected as the optimal parameter. In **Supplementary**

186    **Figure 1**, the objective curve illustrates that the training time was sufficient for model convergence. In

187    this analysis, a shared signature comprising 473 genes between SCZ and BIP was identified, while two

188    disease-specific signatures containing 37 genes for SCZ and 152 genes for BIP, respectively, were found.

189

190

191

192    **Discussion**

193    We here present dsMTL – a secure, federated multi-task learning package for the programming

194    language R, building on DataSHIELD as an ecosystem for privacy-preserving and distributed analysis.

195    Multi-task learning allows the investigation of research questions that are difficult to address using

196    conventionalML, such as the identification of heterogeneous, albeit related, signatures across datasets.

197    The implementation of a privacy-preserving framework for the distributed application of MTL is an

198    essential requirement for the large-scale adoption of MTL. Using such a distributed server setup, we

199    demonstrate the applicability and utility of dsMTL to identify biomarker signatures in different settings.

200    For applications where the target biomarker signatures are different, but relate to an overlapping set

201    of features (explored here as the 'heterogeneous' case), conventional machine learning would not be

202    a meaningful algorithm choice. We show that MTL is able to identify the target signatures with high

203    confidence and may thus be a reasonable choice for a diverse set of interesting analyses. As mentioned

204    above, a particularly noteworthy application is comorbidity modeling, where the target signatures

205     index the shared (although potentially heterogeneously manifested) biology of multiple, clinically

206     comorbid conditions. Such analyses could potentially be a powerful, machine learning-based extension

207     of comorbidity modeling approaches based on univariate statistics that have already been very useful

208     for characterizing the shared biology of comorbid illness[17]. We show that unsupervised MTL can

209     disentangle the shared from cohort-specific effects, demonstrating its potential utility for comorbidity

210     analysis. Other applications for this method include the analysis of biological patterns shared across

211     clinical symptom domains, between clinical and demographic characteristics, or with digital measures,

212     such as ecological momentary assessments.

213     The use of dsMTL follows the concept of the so-called "freely composing script" in the DataSHIELD

214     ecosystem. It organizes a given dsMTL workflow as a free composition of dsMTL, DataSHIELD, and local

215     R commands (e.g. R base functions, customer-defined functions and CRAN packages) into a script, such

216     that the geo-distribution of datasets and the federated computation are transparent to users. This

217     concept is similar to that of the "freely composing apps" used in a recently presented federated ML

218     application[18], which allows flexible scheduling of functions in the form of apps and improves the

219     federated data analysis flexibility for users. In addition to dsMTL, other packages in the DataSHIELD

220     ecosystem exist for e.g. "big data" storage and management[19], various statistical tests[7,19] and deep

221     learning[19,20].

222     Interesting future developments of the dsMTL approach could include the implementation of

223     asynchronous communication, which provides a probabilistically approximate solution but faster

224     convergence[21,22]. Furthermore, integration of other popular systems for ML, such as tensorflow[23], for

225     which interfaces with the R language already exist, would provide valuable additions to the DataSHIELD

226     system. Finally, a noteworthy consideration is an architecture underlying the distributed data

227     infrastructure. DataSHIELD builds on a centralized ("client-server") architecture and each data provider

228     needs to install a well-configured data warehouse. Such infrastructure is suitable for long-term

229     collaboration scenarios and large consortia projects that conduct a broad spectrum of complex

230     analyses requiring high flexibility. However, in other scenarios that require more temporary and easy-

231    compute collaboration setups, a server-free or decentralized architecture[24] might be more suitable,

232    because the cost of data provider for participating is low.

233    In conclusion, the dsMTL library for the programming language R provides an easy-to-use framework

234    for privacy-preserving, federated analysis of geographically distributed datasets. Due to its ability to

235    disentangle shared and cohort-specific effects across these datasets, dsMTL has numerous interesting

236    application areas, including comorbidity modeling and translational research focused on the

237    simultaneous prediction of different outcomes across datasets.

238

239

240    **Methods**

241    *Modeling*

242    All methods part of dsMTL share the identical form,

243
$$\min_{\theta} \mathcal{L}(\theta) + \lambda S(\theta) + C\aleph(\theta)$$

244    where $\mathcal{L}(\theta)$ is the data fitting term (or loss function), the major determinant of the solutions obtained

245    from model training. $\aleph(\theta)$ and $S(\theta)$ are the penalties of $\theta$ with the aim to incorporate the prior

246    information. $\aleph(\theta)$ is a non-smooth function and able to create sparsity, while $S(\theta)$ is smooth. $\lambda$ and $C$

247    are the hyper-parameters to control the strength of the penalties. More technical details can be found

248    in the supplementary methods.

249    In dsMTL, two approaches for sharing information across cohorts are included, 1) shared parameters

250    and 2) cross-task regularization, leading to a slightly different distributed computation. The shared

251    parameters are estimated using all cohorts. For cross-task regularization, the cohort-specific

252    parameters are estimated using only the local data, and then tuned by considering parameters from

253    other cohorts.

254     *Efficiency*

255     Most dsMTL methods aim at training an entire regularization tree. The determination of the $\lambda$

256     sequence controls the tree's growth and is essential for computational speed. The $\lambda$ sequence should

257     be accurately scaled to both capture the highest posterior and avoid overwhelming computations.

258     Inspired by a previous study[25], we estimate the largest and smallest $\lambda$ from the data by characterizing

259     the optima of the objective using the first-order optimal condition and then interpolate the entire $\lambda$

260     sequence on a log scale (see supplementary methods for more details). In addition, several options are

261     provided to improve the speed of the algorithms by decreasing the precision of the results, i.e., 1) the

262     number of digits of parameters for transformation can be specified to reduce the network latency; 2)

263     several termination rules are provided, some of which are relaxed; 3) the depth of the regularization

264     tree can be shortened. More details can be found in supplementary methods.

265     Besides the efficiency of the federated ML/MTL methodology, the import/export of "big data" cohorts

266     is also crucial for computational efficiency, where e.g. uncompressed GWAS data requires tens of

267     gigabytes, leading to time-consuming data import. dsMTL was designed to support a wide variety of

268     data types. For this, an architecture package resourcer[19] developed by the DataSHIELD community was

269     incorporated to facilitate the efficient import and export of large-scale datasets in compressed formats.

270     For example, in DataSHIELD, GWAS data of the PLINK file formats can be read and processed using the

271     software PLINK[26] as the backend[19].

272     *Security*

273     dsMTL was developed based on DataSHIELD[8], which provides comprehensive security mechanisms not

274     specific to machine learning applications. For example, 1) DataSHIELD requires the data analysis to

275     only occur behind the firewall; 2) each server is only allowed to communicate with a set of clients with

276     fixed IP addresses; 3) the network communication is protected by an SSL protocol; 4) an R parser[8]

277     implemented on the server rejects the calling of unwanted functions; and 5) the so-called 'disclosure

278     control'[8] on the server ensures that the returned response does not contain any disclosive information.

279    In addition, several permissions can be set by the data providers to fully control the usage of their data.

280    These permissions describe the degree of accessibility of data and functions on the server i.e. *"which*

281    *users* can perform *what actions* on *what data"*. In an extremely secure example, a user could be

282    granted to check the summary of a given dataset but cannot perform any actions because no functions

283    were granted. With these settings, DataSHIELD allows customizing the security protection strategies

284    according to the specific requirements of the applications. For statistical and machine learning analyses,

285    DataSHIELD assumes that summary statistics are safe to share.

286    dsMTL inherits all these security mechanisms. In addition, we considered potential ML-specific privacy

287    leaks, such as membership inference attacks[27] and model inverse attacks[28]. Inverse attacks aim at

288    extracting the individual observation-level information from the models. Membership inference

289    attempts to decide if an individual was included in a given training set using the model. All these

290    techniques require a complete model for inference. Since multi-task learning returns multiple matrices,

291    returning an incomplete model could be one strategy against these attacks. For example, dsMTL_iNMF

292    in dsMTL only returns the homogenous matrix (H), whereas the cohort-specific components $(V_k, W_k)$

293    never leave the server. For example, in a two-server scenario, one (H) out of five output matrices is

294    transmitted between the client and the servers. With such an incomplete model, inverse construction

295    of the raw data matrix becomes difficult, and the risk of an inverse attack and membership inference

296    is reduced. For most biomedical analyses, the H matrix is sufficient for subsequent studies. In addition,

297    if the analyst was authorized to access the raw data of the server, the so-called "data key mechanism"

298    (see supplement) would allow the analyst to retrieve all component matrices. For supervised multi-

299    task learning methods in dsMTL, all models have to be aggregated within the clients, and thus we

300    suggest the data providers enable the option on the server that rejects a returned coefficient vector

301    containing parameter numbers exceeding the number of subjects. In this way, the model is not

302    saturated and more robust to an inverse attack.

303    *Proof of concept with simulation and actual data*

304 Two case studies and speed-tests were conducted to demonstrate the suitability of dsMTL methods to

305 analyze heterogeneous cohorts, compared to federated ML methods and ensemble of local models

306 regarding the prediction performance, interpretability and computational speed. An overview of

307 methodological aspects related to the case studies is detailed below. For an extensive methodological

308 description, please see the supplementary Methods.

309 **Case study 1.** In this case study, the heterogeneous cohorts were generated with the same set of

310 outcome-associated genes. These however showed different directionality of their respective

311 associations with the outcome. A three-server scenario was simulated. 150 out of 500 features with

312 random signs across cohorts were simulated. Seven tests were created for simulating different n/p

313 ($\frac{\text{sample size}}{\text{gene number}}$) ratios. The n/p ratio was $\{1.2, 1, 0.9, 0.6, 0.5, 0.3, 0.15\}$ with the number of subjects

314 $\{600, 500, 450, 300, 250, 150, 75\}$ for each test. 500 genes were created for each server. The test

315 sample consisted of 200 subjects for each server. Data were generated as follows:

316 Given gene number $p = 500$, the models of three cohorts were $\{w^{(1)}, w^{(3)}, w^{(3)}\}$ where $w^{(.)} = p \times 1$.

317 A shared signature comprising 150 genes was generated for each $w^{(.)}$ but with random signs, $w^{(.)}_i =$

318 $\begin{cases} 2 \times (\rho - 0.5) \times N(1, 0.1) & 1 < i < 150 \\ 0 & \text{others} \end{cases}$, $\rho \sim \text{Bernoulli}(\frac{1}{2})$. The expression values of each subject

319 across cohorts were generated as $x = 1 \times p$ where $x_j \sim N(0,1)$. The numeric outcome (e.g. symptom

320 severity) $y = xw^{(i)}$ in cohort $i$ was standardized in a normal distribution $N(0, 1)$, then model-

321 irrelevant noise with 50% of the variance of the true signal was added $y = y + N(0, 0.5)$.

322 dsMTL_L21 and dsLasso were trained as the federated learning system, and the hyper-parameter was

323 selected using 10 fold in-cohort cross-validation. For glmnet, the ensemble technique was only applied

324 on the gene selection due to the consistent gene set of their signatures. The mean squared error (mse)

325 was used as the measure of prediction performance. To account for the sampling variance, we

326 repeated each analysis 100 times.

327   **Case study 2.** In this case study, two heterogeneous RNA-seq cohorts were created to simulate a

328   comorbidity analysis, where the genes were separated to be part of either a shared signature among

329   cohorts, cohort-specific signatures or diagnosis-unassociated genes. The dsMTL_iNMF was compared

330   to the ensemble of local NMF regarding the selection accuracy of shared/cohot-specific genes, in

331   particular impacted by the severity of heterogeneity. Here the severity of heterogeneity refers to the

332   proportion of the genes harbored by the shared signature over all diagnosis-associated genes. The data

333   simulation protocol for RNA-seq data can be found in the **Supplementary Methods**.

334   A two-server scenario was simulated. As shown in **Supplementary Table 1**, for the data of each server,

335   1000 genes and 200 subjects were simulated, 50% of the genes were diagnosis-unassociated and the

336   remaining genes were part of the disease signature. The genes comprised by shared signatures were

337   identical for data of two servers, and the genes comprised by cohort-specific signatures did not overlap.

338   The case-control ratio was balanced for each server. Four tests were performed by varying the

339   proportion of genes in the shared signature over all diagnosis-associated genes from 20% to 80%.

340   The training of dsMTL_iNMF results in three outputs related to the original input data: the shared gene

341   'exposure' (H), cohort-specific gene 'exposure' (V) and sample 'exposure' (W). We measured the

342   association between the sample exposure and the diagnosis as the weight of each latent factor. The

343   shared( or specific) gene signature was identified as the weighted summation of the shared (or specific)

344   gene exposures over latent factors. To quantify the important genes related to a given signature, we

345   binarized the gene signature according to the mean (0-1 vector, values larger than the mean were

346   assigned). To assess the performance of the gene identification, we associated the selected genes set

347   with the ground truth (0-1 vector, signature genes were 1). The assessment was applied to shared and

348   cohort-specific genes in parallel. Based on this metric, three gene sets were derived as output from

349   dsMTL_iNMF, called dsMTL_iNMF-H, dsMTL_iNMF-V1 and dsMTL_iNMF-V2, and these related to the

350   shared, cohort 1 specific and cohort 2 specific gene signature, respectively. The same strategy was

351   applied to analyze the ensemble of local NMF models. For each cohort, the specific gene signature was

352   the weighted summation of gene exposure over latent factors, and then binarized as the specific gene

353    set (called local-NMF1 and local-NMF2). The shared gene signature was identified as the sum of the

354    specific gene signature over cohorts, and then binarized as the shared gene set(NMF-bagging). We

355    then compared 1) NMF-bagging and dsMTL_iNMF-H for the accuracy related to the isolation of shared

356    genes; 2) dsMTL_iNMF-V1 and local-NMF1 as well as dsMTL_iNMF-V2 and local-NMF2 for the accuracy

357    of isolating cohort-specific genes.

358    **Computational speed of supervised dsMTL.** We aimed at identifying the efficiency of supervised

359    dsMTL using real molecular data and given the real network latency. Four independent schizophrenia

360    case-control cohorts were used for this analysis. The training cohorts consisted of three datasets

361    comprising prefrontal cortex gene expression data (available from the GEO repository under accession

362    numbers GSE53987, GSE21138 and GSE35977). A detailed description of these datasets can be found

363    in their respective original publications[29-31]. The dataset used for algorithm testing was from the HBCC

364    (n=422) cohort comprising genome-wide gene expression data quantified by microarray (dbGAP ID:

365    phs000979.v3.p2). A detailed description of this dataset can be found in the original publication[32]. As

366    shown in **Supplementary Table 2**, three servers were used for training algorithms. Two servers were

367    held at the Central Institute of Mental Health, Mannheim while the third was positioned at the

368    BioQuant institute, Heidelberg.

369    Using this data, we repeated a previously described analysis[13], in order to evaluate computational

370    speed in a federated analysis setting. Here we show the formulation of the mean regularized MTL using

371    dsMTL_net:

372    The cohort-level batch effect was assumed to be Gaussian noise affecting the true coefficient of gene

373    i and cohort j $w_{ij} = w_i + \epsilon_j$, $\epsilon_j \in N(\mu, \sigma)$. Hence, the average model $\overline{w_i}$ across cohorts was an

374    unbiased estimator for the true coefficient, and therefore the squared penalty $\left|w_{ij} - \overline{w_i}\right|^2$ was

375    incorporated to penalize the departure of each model j to the mean. The complete formulation was

376    $$\min_W \sum_{k=1}^{3} \sum_{i=1}^{n_k} \frac{1}{n_k} \log(1 + e^{-Y_i^{(k)}\left(X_i^{(k)}W_{,k}\right)}) + \lambda||W||_1 + C||WG||_2^2,$$

377        where $G = \begin{bmatrix} \frac{2}{3} & 0 & \frac{-1}{3} & \frac{2}{3} & \frac{-1}{3} & 0 \\ \frac{-1}{3} & \frac{2}{3} & 0 & 0 & \frac{2}{3} & \frac{-1}{3} \\ 0 & \frac{-1}{3} & \frac{2}{3} & \frac{-1}{3} & 0 & \frac{2}{3} \end{bmatrix}$
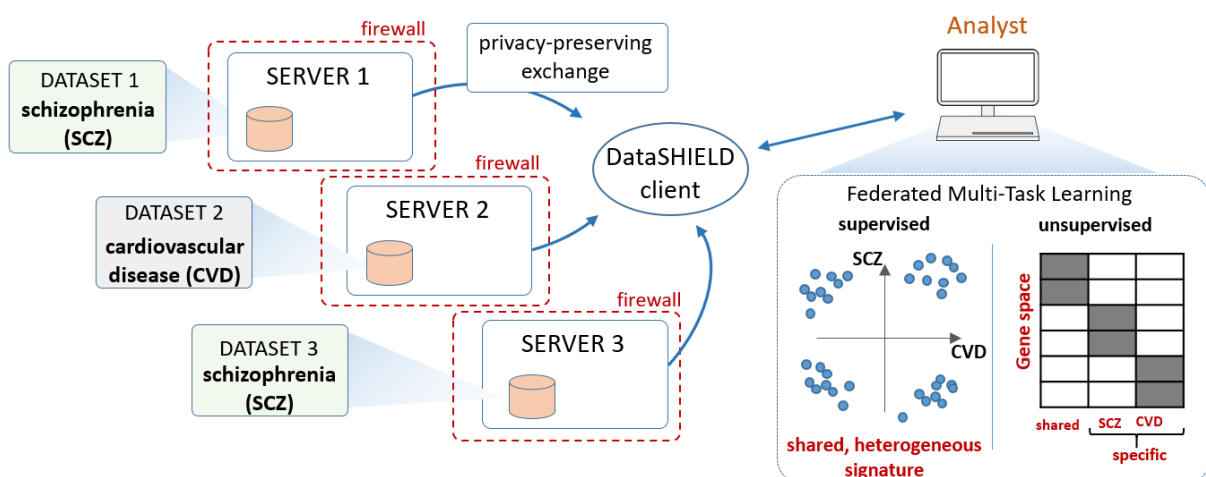
378

379     **Computational speed of unsupervised dsMTL.** Here, we analyzed the time efficiency in applying

380     dsMTL_iNMF on two real datasets based on the real network latency. Two processed RNA-seq case-

381     control cohorts comprising patients with schizophrenia (GSE164376[33] ) and bipolar disorder

382     (GSE134497[34]) were retrieved from the GEO database and converted into a matrix format for the

383     analysis. As shown in **Supplementary Table 4**, the data were stored on servers in Mannheim and

384     Heidelberg.

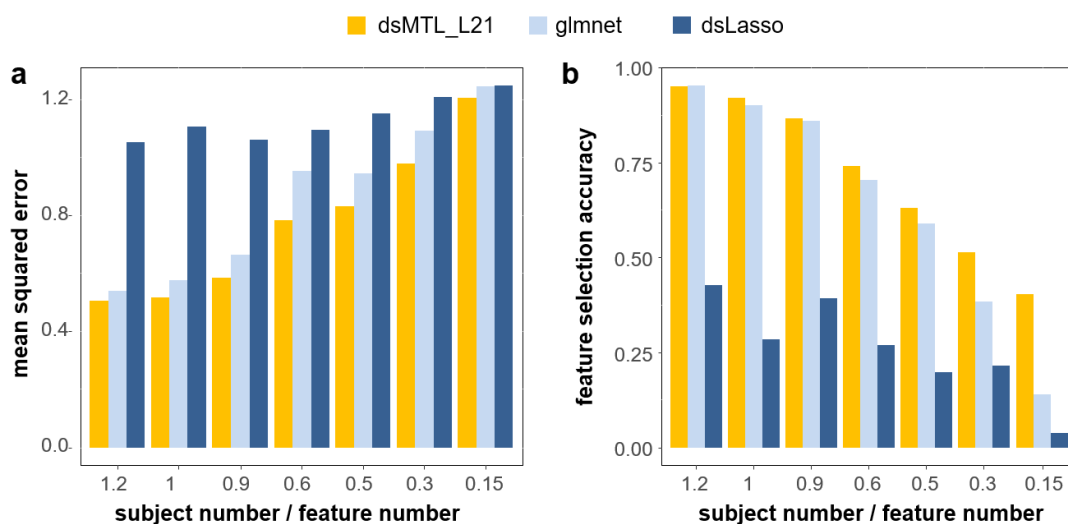385

386

387

388     **Figures**



389

390     **Figure 1. Schematic illustration of dsMTL using comorbidity modeling of schizophrenia and**

391     **cardiovascular disease as an example**. Multiple datasets stored at different institutions are used as a

392     basis for federated MTL. dsMTL was developed in the DataSHIELD ecosystem, which provides

393    functionality regarding data management, transmission and security. Data are analyzed behind a given

394    institution's firewall and only algorithm parameters that do not disclose personally identifiable

395    information are exchanged across the network. dsMTL contains algorithms for supervised and

396    unsupervised multi-task machine learning. The former aims at identifying shared, but potentially

397    heterogeneous signatures across tasks (here, diagnostic classification for schizophrenia and

398    cardiovascular disease). Unsupervised learning separates the original data into shared and cohort-

399    specific components, and aims at revealing the corresponding outcome-associated biological profiles.
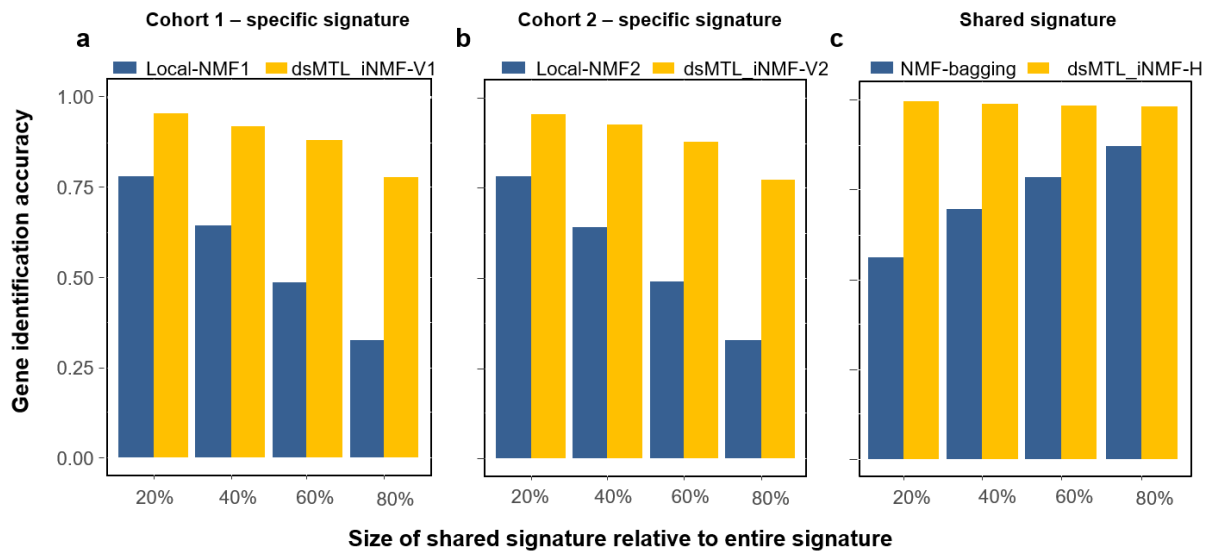
400



401

402    **Figure 2. Analysis of 'heterogeneous' signatures of continuous outcomes in simulated data stored**

403    **on three servers**. The figure shows the **a)** prediction accuracy expressed as the mean squared error

404    and **b)** the feature selection accuracy for different subject/feature number ratios. The respective

405    values were averaged across the three servers, and across 100 repetitions, in order to account for the

406    effect of sampling variability.

407

408

409

410



411

412  **Figure 3. The gene identification accuracy for shared and specific signatures using simulated data. a)**

413  the identification accuracy of important genes for cohort 1. **b)** the identification accuracy of important

414  genes for cohort 2. **c)** the identification accuracy of genes comprised in the shared signature**.** Local-

415  NMF1 and Local-NMF2 were the cohort-specific gene sets identified by local NMF, which were

416  combined into "NMF-bagging" for the shared gene set. dsMTL_iNMF-H was the predicted shared gene

417  set using dsMTL_iNMF. dsMTL_iNMF-V1 and dsMTL_iNMF-V2 were the predicted cohort-specific gene

418  sets identified using dsMTL_iNMF (see Supplementary Figure 1). The proportion of genes harbored by

419  the shared signature was varied from 20% to 80% illustrating the impact of the heterogeneity severity.

420  The model was trained using rank=4 as model parameter. The results for a broader spectrum of rank

421  choices can be found in **Supplementary Figure 2** illustrating that the superior performance of

422  dsMTL_iNMF was not due to the choice of ranks.

423

424

425  **Acknowledgements**

428

429     **Competing interests**

430     AML has received consultant fees from: Boehringer Ingelheim, Elsevier, Brainsway, Lundbeck Int.

431     Neuroscience Foundation, Lundbeck A/S, The Wolfson Foundation, Bloomfield Holding Ltd, Shanghai

432     Research Center for Brain Science, Thieme Verlag, Sage Therapeutics, v Behring Röntgen Stiftung,

433     Fondation FondaMental, Janssen-Cilag GmbH, MedinCell, Brain Mind Institute, Agence Nationale de la

434     Recherche, CISSN (Catania Internat. Summer School of Neuroscience), Daimler und Benz Stiftung,

435     American Association for the Advancement of Science, Servier International. Additionally he has

436     received speaker fees from: Italian Society of Biological Psychiatry, Merz-Stiftung, Forum Werkstatt

437     Karlsruhe, Lundbeck SAS France, BAG Psychiatrie Oberbayern, Klinik für Psychiatrie und

438     Psychotherapie Ingolstadt, med Update GmbH, Society of Biological Psychiatry, Siemens Healthineers,

439     Biotest AG. All other authors have no potential conflicts of interest.

440

441     **References**

442     1.     Jahanshad N, Kochunov PV, Sprooten E, et al. Multi-site genetic analysis of diffusion images
443            and voxelwise heritability analysis: A pilot project of the ENIGMA–DTI working group.
444            *NeuroImage.* 2013;81:455-469.
445     2.     Schizophrenia Working Group of the Psychiatric Genomics C. Biological insights from 108
446            schizophrenia-associated genetic loci. *Nature.* 2014;511(7510):421-427.
447     3.     Kochunov P, Jahanshad N, Sprooten E, et al. Multi-site study of additive genetic effects on
448            fractional anisotropy of cerebral white matter: comparing meta and megaanalytical
449            approaches for data pooling. *NeuroImage.* 2014;95:136-150.
450     4.     Carter KW, Francis RW, Carter K, et al. ViPAR: a software platform for the Virtual Pooling and
451            Analysis of Research Data. *International journal of epidemiology.* 2016;45(2):408-416.
452     5.     Warnat-Herresthal S, Schultze H, Shastry KL, et al. Swarm Learning for decentralized and
453            confidential clinical machine learning. *Nature.* 2021;594(7862):265-270.
454     6.     Zhou J, Liu J, Narayan VA, Ye J, Alzheimer's Disease Neuroimaging I. Modeling disease
455            progression via multi-task learning. *NeuroImage.* 2013;78:233-248.
456     7.     Gaye A, Marcon Y, Isaeva J, et al. DataSHIELD: taking the analysis to the data, not the data to
457            the analysis. *International journal of epidemiology.* 2014;43(6):1929-1944.

8.  Wilson RC, Butters OW, Avraam D, et al. DataSHIELD – New Directions and Dimensions. *Data Science Journal.* 2017;16.

9.  Cao H, Zhou J, Schwarz E. RMTL: An R Library for Multi-Task Learning. *Bioinformatics.* 2018.

10. Yang Z, Michailidis G. A non-negative matrix factorization method for detecting modules in heterogeneous omics multi-modal data. *Bioinformatics.* 2016;32(1):1-8.

11. Xu Q, Xue H, Yang Q. Multi-platform gene-expression mining and marker gene analysis. *International journal of data mining and bioinformatics.* 2011;5(5):485-503.

12. Yuan H, Paskov I, Paskov H, Gonzalez AJ, Leslie CS. Multitask learning improves prediction of cancer drug sensitivity. *Scientific reports.* 2016;6:31619.

13. Cao H, Meyer-Lindenberg A, Schwarz E. Comparative Evaluation of Machine Learning Strategies for Analyzing Big Data in Psychiatry. *International journal of molecular sciences.* 2018;19(11).

14. Zhou J, Yuan L, Liu J, Ye J. A multi-task learning formulation for predicting disease progression. 2011:814.

15. Fujita N, Mizuarai S, Murakami K, Nakai K. Biomarker discovery by integrated joint non-negative matrix factorization and pathway signature analyses. *Scientific reports.* 2018;8(1):9743.

16. Tibshirani R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B (Methodological).* 1996:267-288.

17. Lichtenstein P, Yip BH, Björk C, et al. Common genetic determinants of schizophrenia and bipolar disorder in Swedish families: a population-based study. *The Lancet.* 2009;373(9659):234-239.

18. Matschinske J, Späth J, Nasirigerdeh R, et al. The FeatureCloud AI Store for Federated Learning in Biomedicine and Beyond. *arXiv preprint arXiv:210505734.* 2021.

19. Marcon Y, Bishop T, Avraam D, et al. Orchestrating privacy-protected big data analyses of data from different resources with R and DataSHIELD. *PLoS computational biology.* 2021;17(3):e1008880.

20. Lenz S, Hess M, Binder H. Deep generative models in DataSHIELD. *BMC Med Res Methodol.* 2021;21(1):64.

21. Zhang C, Liu J. Distributed Learning Systems with First-Order Methods. *Foundations and Trends® in Databases.* 2020;9(1):1-100.

22. Xie L, Baytas IM, Lin K, Zhou J. Privacy-Preserving Distributed Multi-Task Learning with Asynchronous Updates. 2017:1195-1204.

23. Dahl M, Mancuso J, Dupis Y, et al. Private machine learning in tensorflow using secure computation. *arXiv preprint arXiv:181008130.* 2018.

24. Warnat-Herresthal S, Schultze H, Shastry KL, et al. Swarm Learning as a privacy-preserving machine learning approach for disease classification. 2020.

25. Friedman J, Hastie T, Tibshirani R. Regularization Paths for Generalized Linear Models via Coordinate Descent. *Journal of Statistical Software.* 2010;33(1).

26. Purcell S, Neale B, Todd-Brown K, et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *American journal of human genetics.* 2007;81(3):559-575.

27. Hu H, Salcic Z, Dobbie G, Zhang X. Membership Inference Attacks on Machine Learning: A Survey. *arXiv preprint arXiv:210307853.* 2021.

28. Fredrikson M, Lantz E, Jha S, Lin S, Page D, Ristenpart T. Privacy in pharmacogenetics: An end-to-end case study of personalized warfarin dosing. Paper presented at: 23rd {USENIX} Security Symposium ({USENIX} Security 14)2014.

29. Lanz TA, Reinhart V, Sheehan MJ, et al. Postmortem transcriptional profiling reveals widespread increase in inflammation in schizophrenia: a comparison of prefrontal cortex, striatum, and hippocampus among matched tetrads of controls with subjects diagnosed with schizophrenia, bipolar or major depressive disorder. *Translational psychiatry.* 2019;9(1):151.

30. Tang B, Capitao C, Dean B, Thomas EA. Differential age- and disease-related effects on the expression of genes related to the arachidonic acid signaling pathway in schizophrenia. *Psychiatry research.* 2012;196(2-3):201-206.

511    31.    Chen C, Cheng L, Grennan K, et al. Two gene co-expression modules differentiate psychotics
512           and controls. *Mol Psychiatry.* 2013;18(12):1308-1314.
513    32.    Fromer M, Roussos P, Sieberts SK, et al. Gene expression elucidates functional impact of
514           polygenic risk for schizophrenia. *Nature neuroscience.* 2016;19(11):1442-1453.
515    33.    A; K, R; K. GSE164376 dataset.
516           https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE164376. Published 2021.
517           Accessed.
518    34.    Kathuria A, Lopez-Lengowski K, Vater M, McPhie D, Cohen BM, Karmacharya R.
519           Transcriptome analysis and functional characterization of cerebral organoids in bipolar
520           disorder. *Genome medicine.* 2020;12(1):34.

521